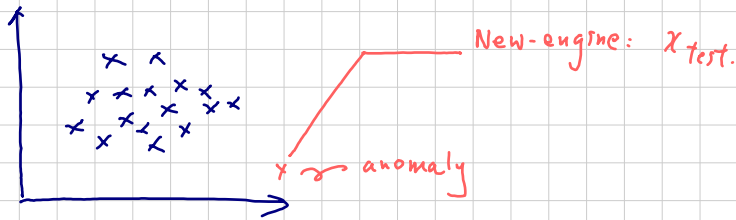


Anomaly Detection

Dataset $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$



Density estimation.

$p(x_{\text{test}}) < \epsilon \rightarrow \text{flag anomaly.}$

$\geq \epsilon \rightarrow \text{ok.}$

how anomaly?

Ex. Fraud detection.

- $x^{(i)}$ = features of user i 's activities
- model $p(x)$ from data.
- identify unusual users by checking $p(x) < \epsilon$.

Ex. Manufacturing

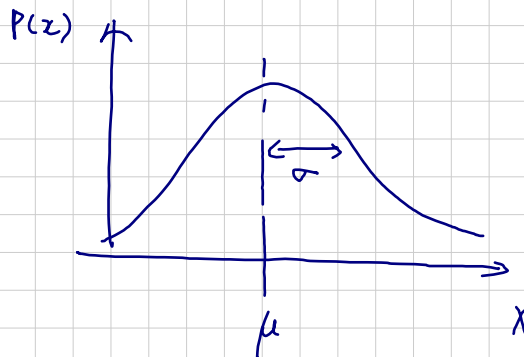
Ex. Monitoring computers in a data centre.

- $x^{(i)}$ = features of machine i .
- x_1 = memory use, x_2 = number of disk accesses / sec.
- x_3 = CPU load, x_4 = CPU load / network traffic.

Gaussian Distribution

$$X \in \mathbb{R}, \quad X \sim \mathcal{N}(\mu, \sigma^2)$$

↗ "normal"
| |
mean variance.

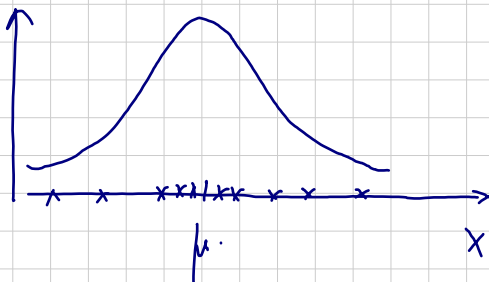


$$p(x; \mu, \sigma^2) \triangleq \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Parameter Estimation

Dataset $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$, $x^{(i)} \in \mathbb{R}$.

If we expect $x^{(i)} \sim \mathcal{N}(\mu, \sigma^2)$



$$\hat{\mu} = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\hat{\sigma}^2 = \frac{1}{m} \sum_{i=1}^m [x^{(i)} - \mu]^2$$

↙
 $n-1$ (sometimes)

Algorithm

Training set $\{x^{(1)}, \dots, x^{(m)}\}$. Each example is $x \in \mathbb{R}^n$.

$$p(x) = \prod_{i=1}^n p(x_i; \mu_i, \sigma_i^2)$$

$$x_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$$

1. Choose features x_i that you think might be indicative of anomalous examples.
2. Fit parameters μ_1, \dots, μ_n , $\sigma_1^2, \dots, \sigma_n^2$ ($j=1, \dots, n$)

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \quad ; \quad \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

$$\underline{\mu} = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix} = \frac{1}{m} \sum_{i=1}^m \underline{x}^{(i)}$$

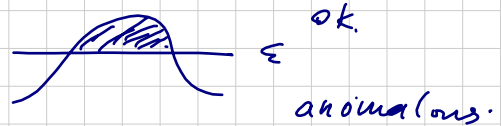
3. Given a new example x

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)$$

$$= \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

$\geq \epsilon \rightarrow \text{normal}$
 $< \epsilon \rightarrow \text{anomaly}$

"Assuming a Gaussian distribution, determine anomaly based on a threshold ϵ "



Developing and Evaluating an Anomaly Detection

• real number evaluation.

Algorithm evaluation

• fit model $p(x)$ on training set. $\{x^{(1)}, \dots, x^{(m)}\}$

• on a CV / test x , predict. $y = \begin{cases} 1, & p(x) < \epsilon \text{ (anomaly)} \\ 0, & p(x) > \epsilon \text{ (normal)} \end{cases}$

↳ skewed data set. (should true/positive, precision/recall, F1 score)

Anomaly Detection vs. Supervised Learning

Examples

$1 < \underline{0}$ (negative examples)
 \downarrow
 $\text{fit } p(x)$

$1, 0$

Many diff. types of anomalies

Enough positive examples for algorithm to get a sense of what positive examples are likely.

Choosing what Features to use

• Non-gaussian $\rightarrow \log. (x+c)$, $x^{\frac{1}{d}}$...
(transformation) \rightarrow constant

• error analysis for anomaly detection

- Want: $p(x)$ large for normal examples x
 $p(x)$ small for anomalous examples x .

- common problem: $p(x)$ is comparable (say, both large) for normal and anomalous examples.

Multivariate Gaussian.

$x \in \mathbb{R}^n$. Don't model $p(x_1), p(x_2), \dots$, etc. separately.

Model $p(x)$ all in one-go.

Parameters: $\mu \in \mathbb{R}^n$, $\Sigma \in \mathbb{R}^{n \times n}$. [covariance matrix].

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} \cdot |\Sigma|^{\frac{1}{2}}} \exp. \left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right)$$

\hookrightarrow determinant.

$$\Sigma = \begin{pmatrix} 0.6 & 0 \\ 0 & 0.6 \end{pmatrix} \rightarrow \text{variance shrunk} \rightarrow \text{narrower distribution}$$

$$\Sigma = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix} \rightarrow \text{flatter.}$$

Anomaly Detection with the multivariate Gaussian.

1) Fit model $p(x)$ by setting

$$\mu = \frac{1}{m}, \quad \Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu) \cdot (x^{(i)} - \mu)^T.$$

2) Given a new example x , compute

$$p(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right).$$

$$\text{anomaly} \Leftrightarrow p(x) < \epsilon.$$

• Original model $p(x) = \prod p(z_i; \mu_i, \sigma_i^2) \Leftrightarrow p(z; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$

where $\Sigma = \begin{pmatrix} \sigma & \emptyset & \emptyset \\ \emptyset & \dots & \emptyset \\ \emptyset & \emptyset & \sigma \end{pmatrix}$

special case.

• When to use?

original model.

$$\prod p(x_i; \mu, \sigma)$$

- Manually create features to capture anomalies where x_1, x_2 take unusual combinations of values. $x_3 = \frac{x_1}{x_2}$

• m can $< n$.

multivariate Gaussian

$$p = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

- Automatically captures correlations between features.

• Computationally expensive.

• Σ invertible $\Leftrightarrow m > n$.

($m > 10n$ in practice)

↳ redundant features when $n \uparrow$

Recommender system (use content of movie to predict ratings)

1. Problem formulation: predict users movie rating \rightarrow generate recommendations

2. Content-based recommendation

• for each movie, form "feature vector" $= \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix} \rightarrow \begin{matrix} 1. \\ \text{romance} \\ \text{action} \end{matrix}$

• for each user j , learn parameter $\theta^{(j)} \in \mathbb{R}^3$. Predict user j as rating movie i with $(\theta^{(j)})^T x^{(i)}$ starts.

3. Formulation.

• $r(i,j) = 1$ if user j rated movie i . (\emptyset otherwise)

• $y^{(i,j)} = \text{rating}$ "

• $\theta^{(j)}$ = parameter vector for user j .

• $x^{(i)}$ = feature vector for movie i .

• For user j , movie i , predictive $(\theta^{(j)})^T x^{(i)}$.

• $m^{(j)} = \#$ movies rated by user j

To learn $\theta^{(j)}$:

$$\min_{\theta^{(j)}} \underbrace{\frac{1}{2m^{(j)}} \sum_{i:r(i,j)=1} [(\theta^{(j)})^T x^{(i)} - y^{(i,j)}]^2}_{\text{data fit}} + \underbrace{\frac{\lambda}{2m^{(j)}} \sum_{k=1}^n (\theta_k^{(j)})^2}_{\text{regularization}}, \quad \theta^{(j)} \in \mathbb{R}^{n+1}.$$

constant, can pull out.

4. Optimization.

$$\textcircled{1} \min_{\theta^{(1)} \dots \theta^{(m)}} \frac{1}{2} \sum_{j=1}^{m_u} \sum_{i:r(i,j)=1} []^2 \rightarrow \frac{\lambda}{2} \sum \sum (\theta_k^{(j)})^2$$

$J(\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(m)})$

② gradient descent.

Feature Learning: Collaborative Filtering

- Given a dataset, but no knowledge of each movie.
- However, has user data θ parameter.

1. Formulation.

Given: $\theta^{(1)}, \dots, \theta^{(n_u)}$ to learn $x^{(i)}$:

$$\min_{x^{(i)}} \frac{1}{2} \sum_{j: r(i,j)=1} [(\theta^{(j)})^T x^{(i)} - y^{(i,j)}]^2 + \frac{\lambda}{2} \sum_{k=1}^n (x_k^{(i)})^2$$

↓ sum over all movies.

$$\min \frac{1}{2} \sum_{i=1}^{n_u} \left(\dots \right)$$

2. Summary

$\left\{ \begin{array}{ll} \text{Given } x^{(1)}, \dots, x^{(n_u)} \text{ [and movie ratings]} & \xrightarrow{\text{est.}} \theta^{(1)} \dots \theta^{(n_u)} \\ \text{Given } \theta^{(1)}, \dots, \theta^{(n_u)} & \rightarrow x^{(1)} \dots x^{(n_u)} \end{array} \right.$
Can cascade to improve.

3. Collaborative Filtering Algorithm

- Minimize over $(i,j): r(i,j)=1$.
- $x \in \mathbb{R}^n$, $\theta \in \mathbb{R}^n \rightarrow$ No need to hardcode $x_0=1$, since we are generating the features.
- 1. Initialize $x^{(1)}, \dots, x^{(n_u)}$; $\theta^{(1)}, \dots, \theta^{(n_u)}$ to small random values.
- 2. Min. $J(x^{(1)}, \dots, x^{(n_u)}, \theta^{(1)}, \dots, \theta^{(n_u)})$ using gradient descent.
- 3. For a user with parameters θ , and a movie w/ learned features x , predict a star rating of $\theta^T x$.

Vectorization: Low Rank Matrix Factorization

$$Y = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} \rightarrow y^{(i,j)} \text{ movie } i, \text{ user } j$$

Predicted ratings: $[\theta^{(j)T} \cdot x^{(i)}]$, $j=1 \dots n_u$
 $i=1 \dots n_m$

$$X = \begin{pmatrix} -x^{(1)T} - \\ -x^{(2)T} - \\ \vdots \\ -x^{(n_m)T} - \end{pmatrix} \rightarrow \text{each movie example in row.}$$

$$\Theta = \begin{pmatrix} -\theta^{(1)T} - \\ -\theta^{(2)T} - \\ \vdots \\ -\theta^{(n_u)T} - \end{pmatrix}$$

$$\boxed{Y = X \Theta^T} \rightarrow \text{low rank matrix factorization.}$$

Finding Related Movies

For each product i , we learn a feature vector $x^{(i)} \in \mathbb{R}^n$

How to find movies j related to movie i ?

might not be human understandable.

small $\|x^{(i)} - x^{(j)}\| \rightarrow$ two movies are "similar", share similar viewers

Mean Normalization

New users: $\min_{\theta} \underbrace{0} + \frac{\lambda}{2} \sum \sum (x_k^{(i)})^2 + \frac{\lambda}{2} \sum \sum (\theta_k^{(i)})^2$

\downarrow
 $\theta = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$

Mean normalization \rightarrow subtract average of movie

\rightarrow normalise each row to μ .
course

\downarrow
learn $\theta^{(i)}, x^{(i)}$.

For user j , movie i . predict:

$$\underbrace{(\theta^j)^T \cdot (x^{(i)})}_0 + \underline{\underline{\mu_i}}$$